1    This application is submitted in the name of the following inventor(s):

2

3    <u>Inventor</u>                    <u>Citizenship</u>                    <u>Residence City and State</u>

4    English, Robert M.            United States                  Menlo Park, California

5

6         The assignee is <u>Network Appliance, Inc.</u>, a corporation having an office at

7    495 East Java Drive, Sunnyvale, California, 94089.

8

9                              <u>Title of the Invention</u>

10

11                    Low-Overhead Threads in A High-Concurrency System

12

13                            <u>Background of the Invention</u>

14         This application claims the benefit of U.S. Provisional Application No.

15    60/195,732, filed 4/7/00 (Attorney Docket number 103.1032.01).

16

17    *1.    Field of the Invention*

18

19         This invention relates to low-overhead threads in a high-concurrency sys-

20    tem, such as for a networked cache or file server.

21

1   *2.    Related Art*

2

3          In many computing systems, it is desirable in certain circumstances to be

4   able to process, relatively simultaneously (such as in parallel), a relatively large number

5   of similar tasks.  For example, the same or similar tasks could be performed by a server

6   device (such as a file server) in response to requests by a number of client devices.  One

7   such circumstance is in a networked cache or file server, which maintains and processes a

8   relatively large number of sequences of requests (sometimes called "connections"), so as

9   to couple an information requester (such as a web client) to one or more information pro-

10   viders, which are also coupled to the same internetworking system.  One known method

11   in which an individual processor or a multiprocessor system is able to maintain a high de-

12   gree of concurrency is for the system to process each connection using a separate proc-

13   essing thread.  A "thread" is a locus of control within a process, indicating a spot within

14   that process that the processor is then currently executing.  In general, a thread has a rela-

15   tively small amount of state information associated therewith, generally consisting only of

16   a calling stack and a relatively small number of local variables.

17

18          High concurrency systems, such as networked caches and file servers used

19   in an internetworking system, must generally maintain a large number of threads.  Each

20   information requester has its own separate connection for which the network cache or file

21   server must maintain some amount of state information.  Each such separate connection

22   requires only a small amount of state information, such as approximately 100 to 200 bytes

1    of information. Since there are in many cases a relatively large number of individual

2    connections, it would be desirable to be able to maintain state information about each

3    such connection using only a relatively minimal amount of memory and processor over-

4    head, while simultaneously maintaining both relatively reliable programmability and rela-

5    tively high processing speed.

6

7              One problem with known systems is that allocation of state information for

8    individual threads does not generally scale well. One of the problems with relatively

9    large numbers of individual threads is that of allocating memory space for a calling stack

10   for each one of those threads. In a first set of known systems, stack space for individual

11   threads is allocated statically; this has the drawback that relatively large numbers of

12   threads require a relatively large amount of memory to maintain all such stack spaces.

13   Although the amount of stack space statically allocated for each individual thread can be

14   reduced significantly, this has the drawback that operations that can be performed by each

15   individual thread are similarly significantly restricted. In a second set of known systems,

16   stack space for individual threads is allocated dynamically; this has the drawback that the

17   minimum size for dynamic allocation of memory is generally measured in kilobytes, re-

18   sulting in substantial unnecessary memory overhead. Although virtual memory can be

19   used to store and retrieve stack space for individual threads in smaller increments, this has

20   the drawback that compression and decompression of stack space for individual threads

21   imposes substantial unnecessary processor overhead. In a third set of known systems,

22   such as those using the Java programming language, dynamic memory allocation is used

1 to store and retrieve stack space for individual threads; this has the drawback that each

2 procedure call within each thread imposes substantial unnecessary processor overhead.

3

4 An additional problem is introduced by the particular use made of multi-

5 threading by the WAFL file system (as described in the Incorporated Disclosures). In the

6 WAFL file system, the C language "setjmp" and "longjmp" routines are combined with

7 message passing among threads so as to support high concurrency using threads. In par-

8 ticular, the requester of an initial file request to the WAFL file system packages the re-

9 quest in a message, which the WAFL file system processes using ordinary procedural

10 program code, so long as data is available for processing the request and the thread need

11 not have its execution suspended. If the thread is suspended for any reason (such as if a

12 resource is not available,) the WAFL file system: (1) requests the needed resource, (2)

13 queues the message for signaling when the resource is available, and (3) calls the C rout-

14 ing "longjmp" to return to the origin of the routine for processing the message. Thus, the

15 WAFL file system restarts processing the entire message from the very beginning until all

16 needed resources are available and processing can complete without suspension. While

17 this use of multithreading by the WAFL file system has the advantage that programmers

18 do not need to encode program state when a routine is suspended, it has the disadvantage,

19 when combined with multithreading, that all necessary data structures (to process any ar-

20 bitrary message) must be collected before the entire message can be processed. In an in-

21 ternetworking environment, collecting all such structures can be difficult and subject to

22 error.

1

2      Accordingly, it would be advantageous to provide a technique for creating

3   and using relatively low-overhead threads in a high-concurrency system, such as for a

4   networked cache or file server, that is not subject to drawbacks of the known art.

5

6                          <u>Summary of the Invention</u>

7

8      The invention provides a method and system for providing the functionality

9   of dynamically-allocated threads in a multithreaded system in which the operating system

10  provides only statically-allocated threads.  With this functionality, a relatively large num-

11  ber of threads can be maintained without a relatively large amount of overhead (either in

12  memory or processor time), and it remains possible to produce program code without un-

13  due complexity.

14

15     In a preferred embodiment, a plurality of dynamically-allocated threads are

16  simulated using a single statically-allocated thread, but with state information regarding

17  each dynamically-allocated thread maintained within the single statically-allocated thread.

18  The single statically-allocated thread includes, for each procedure call that would other-

19  wise introduce a new dynamically-allocated thread, a memory block including: (1) a rela-

20  tively small procedure call stack for the new dynamically-allocated thread, and (2) a rela-

21  tively small collection of local variables and other state information for the new dynami-

22  cally-allocated thread. When using multithreading in the WAFL file system, high

1 concurrency among threads can be maintained without any particular requirement that the

2 program code maintain a substantial amount of state information regarding each dynami-

3 cally-allocated thread. Each routine in the WAFL file system that expects to be sus-

4 pended or interrupted need maintain only a collection of entry points into which the rou-

5 tine is re-entered when the suspension or interruption is completed. A feature of the C

6 language preprocessor allows the programmer to generate each of these entry points

7 without substantial additional programming work, with the aid of one or more program-

8 ming macros.

9

10 The invention provides an enabling technology for a wide variety of appli-

11 cations for multithreaded systems so as to obtain substantial advantages and capabilities

12 that are novel and non-obvious in view of the known art. Examples described below pri-

13 marily relate to networked caches and file servers, but the invention is broadly applicable

14 to many different types of automated software systems.

15

16 Brief Description of the Drawings

17

18 Figure 1 shows a block diagram of a system for providing functionality of

19 low-overhead threads in a high-concurrency system, such as for a networked cache or file

20 server.

21

1         Figure 2 shows a process flow diagram of a system for providing function-

2   ality of low-overhead threads in a high-concurrency system, such as for a networked

3   cache or file server.

4

5                   <u>Detailed Description of the Preferred Embodiment</u>

6

7         In the following description, a preferred embodiment of the invention is de-

8   scribed with regard to preferred process steps and data structures. Embodiments of the

9   invention can be implemented using general-purpose processors or special purpose proc-

10   essors operating under program control, or other circuits, adapted to particular process

11   steps and data structures described herein. Implementation of the process steps and data

12   structures described herein would not require undue experimentation or further invention.

13

14   *Lexicography*

15

16         The following terms refer or relate to aspects of the invention as described

17   below. The descriptions of general meanings of these terms are not intended to be limit-

18   ing, only illustrative.

19

20   o     client and server — In general, these terms refer to a relationship between two

21          devices, particularly to their relationship as client and server, not necessarily to any

22          particular physical devices.

For example, but without limitation, a particular client device in a first relationship with a first server device, can serve as a server device in a second relationship with a second client device. In a preferred embodiment, there are generally a relatively small number of server devices servicing a relatively larger number of client devices.

o   client device and server device — In general, these terms refer to devices taking on the role of a client device or a server device in a client-server relationship (such as an HTTP web client and web server). There is no particular requirement that any client devices or server devices must be individual physical devices. They can each be a single device, a set of cooperating devices, a portion of a device, or some combination thereof.

For example, but without limitation, the client device and the server device in a client-server relation can actually be the same physical device, with a first set of software elements serving to perform client functions and a second set of software elements serving to perform server functions

As noted above, these descriptions of general meanings of these terms are not intended to be limiting, only illustrative. Other and further applications of the invention, including extensions of these terms and concepts, would be clear to those of ordinary

1    skill in the art after perusing this application.  These other and further applications are

2    part of the scope and spirit of the invention, and would be clear to those of ordinary skill

3    in the art, without further invention or undue experimentation.

4

5    *System Elements*

6

7            Figure 1 shows a block diagram of a system for providing functionality of

8    low-overhead threads in a high-concurrency system, such as for a networked cache or file

9    server.

10

11            A system 100 includes a networked cache or file server (or other device)

12    110, a sequence of input request messages 120, and a set of software elements 130.

13

14            The networked cache or file server (or other device) 110 includes a com-

15    puter having a processor, program and data memory, mass storage, a presentation ele-

16    ment, and an input element, and is coupled to a communication network.  As used herein,

17    the term "computer" is intended in its broadest sense, and includes any device having a

18    programmable processor or otherwise falling within the generalized Turing machine

19    paradigm.  The mass storage can include any device for storing relatively large amounts

20    of information, such as magnetic disks or tapes, optical devices, magneto-optical devices,

21    or other types of mass storage.

22

1    The input request messages 120 include a set of messages requesting the

2    networked cache or file server 110 to perform actions in response thereto. In a preferred

3    embodiment, the actions to be performed by the networked cache or file server 110 will

4    involve access to the mass storage or to the communication network. In a preferred em-

5    bodiment, the input request messages 120 are formatted in a known request protocol, such

6    as NFS, CIFS, HTTP (or variants thereof), but there is no particular requirement for the

7    input request messages 120 to use these known request protocols or any other known re-

8    quest protocols. In a preferred embodiment, the networked cache or file server 110 re-

9    sponds to the input request messages 120 with both: (1) a condign set of responsive ac-

10   tions involving the mass storage or the vacation network, and (2) a condign response to

11   the input request messages 120, the response to the input request messages 120 preferably

12   taking the form of a set of response messages (not shown.)


14   The software elements 130 include a set of programmed routines to be per-

15   formed by the networked cache or file server 110, using the functionality of low-overhead

16   threads and high-concurrency as described herein. Although particular program code is

17   described herein with regard to the programmed routines, there is no particular reason that

18   the software elements 130 must use the specific program code described herein, or any

19   other specific program code.

20

1 *Method of Operation*

2

3    Figure 2 shows a process flow diagram of a system for providing function-

4 ality of low-overhead threads in a high-concurrency system, such as for a networked

5 cache or file server.

6

7    A method 200 includes a set of flow points and a set of steps. The system

8 100 performs the method 200. Although the method 200 is described serially, the steps of

9 the method 200 can be performed by separate elements in conjunction or in parallel,

10 whether asynchronously, in a pipelined manner, or otherwise. There is no particular re-

11 quirement that the method 200 be performed in the same order in which this description

12 lists the steps, except where so indicated.

13

14    At a flow point 210, the networked cache or file server 110 is ready to re-

15 ceive and respond to the input request messages 120.

16

17    At a step 211, the networked cache or file server 110 receives an input re-

18 quest message 120, and forwards that input request message 120 to an appropriate soft-

19 ware element 130 for processing. In a preferred embodiment, the step 211 includes per-

20 forming a calling sequence for the software element 130, including possibly creating a

21 simulated dynamically allocated thread (that is, a thread simulated so as to appear to be

22 dynamically-allocated, hereinafter sometimes called a "simulated thread" or an "S-

1    thread") within which the software element 130 is performed. Thus, the software element

2    130 can be created using program code that assumes that the software element 130 is per-

3    formed by a separate thread and does not demand relatively excessive resources (either

4    memory or processor time.)

5

6    As part of step 211, the networked cache or file server 110 allocates a pro-

7    cedure call block 131 and a local variable block 132, for use by the simulated dynami-

8    cally-allocated thread performed by the software element 130. The procedure call block

9    131 includes a set of input variables for input to the software element 130, a set of output

10   variables for output from the software element 130, and such other stack element as is

11   known in the art of calling stacks for procedure calls. The local variable block 132 in-

12   cludes a set of locations in which to store local variables for the software element 130.

13

14   As part of step 211, the networked cache or file server 110 determines

15   whether the software element 130 is a subroutine of a previously called software element

16   130 in the same simulated thread. If so, the networked cache or file server 110 indicates

17   that fact in a block header 133 for the software element 130, so as to point back to the

18   particular software element 130 that was the parent (calling) software element 130. If

19   not, the networked cache or file server 110 does not indicate that fact in the block call or

20   block header for the software element 130.

21

As part of this step, the networked cache or file server 110 determines whether the software element 130 is to be performed by a new simulated thread. If so, the networked cache or file server 110 adds the new thread block 134 to a linked list 135 of thread blocks 134 to be performed in turn according to a scheduler. In a preferred embodiment, the scheduler simply performs each simulated thread corresponding to the next thread block 134 in round-robin sequence, so that each simulated thread corresponding to a thread block 134 is performed in its turn, until it is suspended or completes. However, in alternative embodiments, the scheduler may select simulated threads in other than a round-robin sequence, so as to achieve a desired measure of quality of service, or other administrative goals.

At a step 212, the networked cache or file server 110 chooses the simulated thread for execution. The simulated thread, with appropriate data completed for the procedure call block 131 and local variable block 132, is performed in its turn, until it is suspended or completes. If the simulated thread is capable of completing its operation without being suspended or interrupted, the scheduler selects the next thread block 134 in the linked list of thread blocks 134 to be performed in turn.

After this step, the method 200 has performed one round of receiving and responding to input request messages 120, and is ready to perform another such round so as to continuously receive and respond to input request messages 120.

1    The method 200 is performed one or more times starting from the flow

2    point 210 and continuing therefrom. In a preferred embodiment, the networked cache or

3    file server 110 repeatedly performs the method 200, starting from the flow point 210 and

4    continuing therefrom, so as to receive and respond to input request messages 120 periodi-

5    cally and continuously.

6

7    *Program Structures*

8

9    A set of program structures in a system for providing functionality of low-

10   overhead threads in a high-concurrency system, such as for a networked cache or file

11   server, includes one or more of, or some combination of, the following:

13   ○    A set of program structures for declaring and creating a dynamically-allocated thread

14        in a system in which threads are usually statically-allocated;

16   ```
     typedef struct {

17       // local variables

18       int arg;      // an example, not necessary

19   } function_msg;
     ```

20

21   In the program structure above, the definition for the structure type "func-

22   tion_msg" includes: (1) the local variables for the dynamically-allocated thread, (2) any

1  input arguments to the dynamically-allocated thread, in this case just the one variable

2  "arg", and (3) any output arguments from the dynamically-allocated thread, in this case

3  none.

4

5  o  A set of program structures for denoting program code entry-points for a simulated

6    thread;

```
static void
    function_sthread(sthread_msg *m)
    {
        function_msg * const msg = m->data;

        STHREAD_START_BLOCK (m);
        // executable C code
        STHREAD_RESTART_POINT (m);          // an example
blocking point
        // executable C code
        STHREAD_COND_WAIT (m, cond (m));  // encapsulated
blocking point
        // executable C code
        STHREAD_END_BLOCK;
        free (msg);
    }
```

24    The program structure above includes, in its definition for the function

25  "function_sthread", an initial program statement obtaining access to the local variables

26  for the simulated thread. This is the statement referring to "m -> data".

The program structure above includes a definition for a start-point for the simulated thread. This is the statement "STHREAD_START_BLOCK (m)", which makes use of a macro defined for the name "STHREAD_START_BLOCK".

The program structure above includes a definition for a restart-point for the simulated thread. This is the statement "STHREAD_RESTART_POINT (m)", which makes use of a macro defined for the name "STHREAD_RESTART_POINT".

The program structure above includes a definition for a conditional-wait point (a possible suspension of the simulated thread) for the simulated thread. This is the statement "STHREAD_COND_WAIT(m, cond(m))", which makes use of a macro defined for the name "STHREAD_COND_WAIT".

The program structure above includes, in its definition for the function "function_sthread", a closing program statement for ending the simulated thread. This is the statement "STHREAD_END_BLOCK", which makes use of a macro defined for the name "STHREAD_END_BLOCK". The program structure above also includes a statement for freeing any data structures used by the simulated thread. This is the statement "free(msg)".

1    The macro definitions for "STHREAD_START_BLOCK",

2  "STHREAD_RESTART_POINT", and "STHREAD_END_BLOCK" collectively form

3  a C language "case" statement.

4

5    o  The macro "STHREAD_START_BLOCK" includes the preamble to the

6       "case" statement:

7

8
```
#define STHREAD_START_BLOCK (m) switch (m -> line) { case 0:
```

9

10   o  The macro "STHREAD_RESTART_POINT" includes an intermediate restart

11      point in the "case" statement:

12

13
```
#define STHREAD_RESTART_POINT(m) case __LINE__: m -> line
= __LINE__
```

16    The restart point uses the C preprocessor to generate tags that the switch

17  statement uses as branch points. The C macro __LINE__ substitutes the line number of

18  the file being processed, so a series of restart points generates a series of unique cases

19  within the switch. Setting m -> line to the case just entered means that if the procedure is

20  re-entered the switch statement will branch to the restart point and continue.

21

1    o   The macro "STHREAD_START_BLOCK" includes the close of the "case"

2    statement:

3

```
#define STHREAD_END_BLOCK }
```

4

5

6    Thus, the C preprocessor generates a "case" statement in response to use of

7    these macros, which allows the programmer to easily specify each of the proper restart

8    points of the routine.

9

10    o   A set of program structures for suspending and restarting simulated threads;

11

```
#define          STHREAD_COND_WAIT(m,          c)          \
STHREAD_RESTART_POINT(m); \
        { \ if (c) \
                sthread_suspend(); \
        }
```

17

18    At an individual restart point, the programmer can use the macro

19    "STHREAD_COND_WAIT" to conditionally either wait for an operation to complete,

20    or to suspend and restart the simulated thread while waiting for resources for the opera-

21    tion to complete.

22

23    o   A set of program structures for initiating simulated threads;

o   The macro "STHREAD_INIT" allocates memory for the simulated thread, sets the C preprocessor value __LINE__ to zero, sets the value of "data" to the private stack area of the particular simulated thread, and sets a value for "handler" to a function passed to the macro as an argument.

```
#define STHREAD_INIT(m, msg, handler) \ m = malloc(sizeof(*m)); \
msg = zalloc(sizeof(*msg)); \ m -> line = 0; \ m -> data = msg; \ m ->
handler = handler
```

o   A set of program structures for actually performing the simulated thread;

```
void
    function(int arg)
    {
        function_msg *msg;
        sthread_msg *m;

        STHREAD_INIT(m, msg, function_sthread);
        msg->arg = arg;

        sthread_run(m);
    }
```

The program structure above includes, in its definition for the function "function", program code for creating the data blocks for the simulated thread, and for placing data in those data blocks. These are the statements "STHREAD_INIT(m, msg, function_sthread)" and "msg -> arg = arg", which make use of a macro defined for the name "STHREAD_INIT".

○ A set of program structures for scheduling performance of simulated threads;

```
switch (m->line) {     // a field in sthread_msg
case 0:
     // executable C code
STHREAD_RESTART_POINT(m);
     // executable C code
STHREAD_RESTART_POINT(m);
     // executable C code
}
```

The program structure above includes, in its definition for the function "function", program code for creating the data blocks for the simulated thread, and for placing data in those data blocks. These are the statements "STHREAD_INIT(m, msg, function_sthread)" and "msg -> arg = arg", which make use of a macro defined for the name "STHREAD_INIT".

1    o   A set of program structures for suspending and resuming performance of simulated

2        threads.

3

```
typedef struct sthread_msg {
    int line;
    void *data;
    void (*handler)(sthread_msg *);
}
jmp_buf sthread_env;
void
sthread_run(sthread_msg *m)
{
    if (!setjmp(sthread_env)) {
        m->handler(m);
        free(m);
    }
}
void
sthread_suspend()
{
    longjmp(sthread_env, 0);
}
sthread_msg *suspended_sthread;
int ready;
int
cond(sthread_msg *m)
{
    if (ready)
```

```
            return 1;
        suspended_sthread = m;
        sthread_suspend();
    }


    int
    set_cond()
    {
        ready = 1;
        if (suspended_sthread) {
            sthread_msg *m = suspended_sthread;
            suspended_sthread = 0;
            sthread_run(m);
        }
    }
    // cond() changed
    sthread_run(suspended_sthread);
```

and

- A set of program structures for performing simulated threads in conjunction with the WAFL file system, as shown above.

*Generality of the Invention*

1

2

3    The invention has general applicability to various fields of use, not neces-

4 sarily related to the services described above.  For example, these fields of use can in-

5 clude devices other than file servers.

6

7    Other and further applications of the invention in its most general form, will

8 be clear to those skilled in the art after perusal of this application, and are within the

9 scope and spirit of the invention.

10

*Technical Appendix*

11

12

13    The technical appendix enclosed with this application is hereby incorpo-

14 rated by reference as if fully set forth herein, and forms a part of the disclosure of the in-

15 vention and its preferred embodiments.

16

*Alternative Embodiments*

17

18

19    Although preferred embodiments are disclosed herein, many variations are

20 possible which remain within the concept, scope, and spirit of the invention, and these

21 variations would become clear to those skilled in the art after perusal of this application.